

An Insight into Gill Microbiome of Eastern Mediterranean Wild Fish by Applying Next Generation Sequencing

Peleg Itay^a, Eli Shemesh^a, Maya Ofek-Lalzar^b, Nadav Davidovich^{a,c}, Yael Kroin^a, Shlomi Zrihan^a, Nir Stern^d, Arik Diamant^a, Natascha Wosnick^e Dalit Meron^a, Dan Tchernov^a, Danny Morick^{a,f,#}

^aMorris Kahn Marine Research Station, Department of Marine Biology, Leon H. Charney School of Marine Sciences, University of Haifa, Israel.

^bBioinformatics Services Unit, University of Haifa, Haifa, Israel.

^cIsraeli Veterinary Services, Bet Dagan, 5025001, Israel.

^dNational Institute of Oceanography, Israel Oceanographic and Limnological Research, Haifa, Israel

^ePrograma de Pós-graduação em Zoologia, Universidade Federal do Paraná, Brazil

^fHong Kong Branch of Southern Marine Science and Engineering, Guangdong Laboratory (Guangzhou), Hong Kong, China.

[#]Corresponding Author: Danny Morick (dmorick@univ.haifa.ac.il)

This PDF file includes:

Figures S1-S8

(Table S1 is an additional Excel file)

Supplementary methods

Supplementary figures

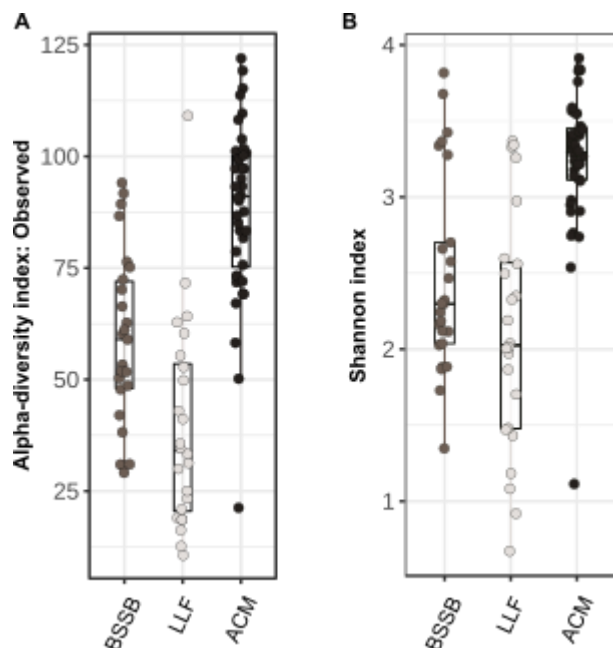


Figure S1: Alpha-diversity analyses comparison between fish

samples using Observed (A) and Shannon index (B). Samples

clustered by fish species.

Observed: Kruskal-Wallis statistic: 44.181, $P < .001$.

Shannon: Kruskal-Wallis statistic: 29.703, $P < .001$.

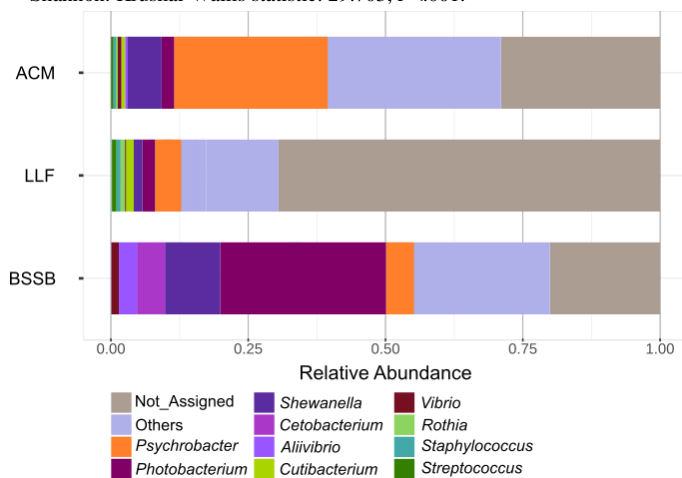


Figure S2: Average relative abundance of the fish species' major

bacterial populations within the gills' community. The graph presents

groups by relative size – from right to left – where the largest (more

abundant) group is on the right and the least is on the left. The grey

colored portions of the bars represent bacteria unidentified to the genus

taxonomic level. Low count bacteria genera were aggregated into the

'Others' group. The three fish species express big differences in prominent

genera, including presence of potential pathogens.

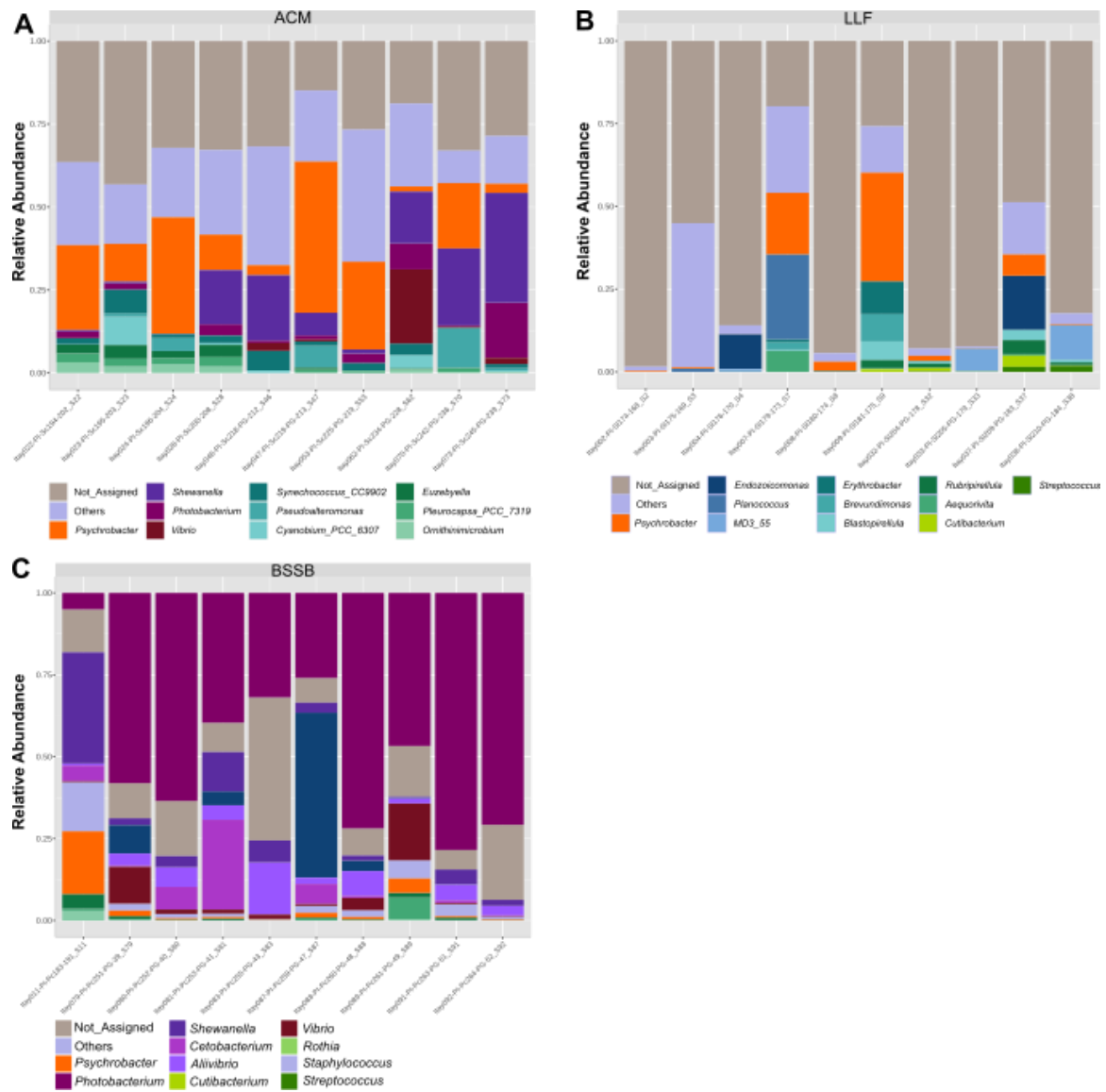


Figure S3: Atlantic Chub Mackerel (ACM; A), Lessepsian Lizardfish (LLF; B) and Bluespotted Seabream (BSSB; C) relative abundance – 10 random samples (each).

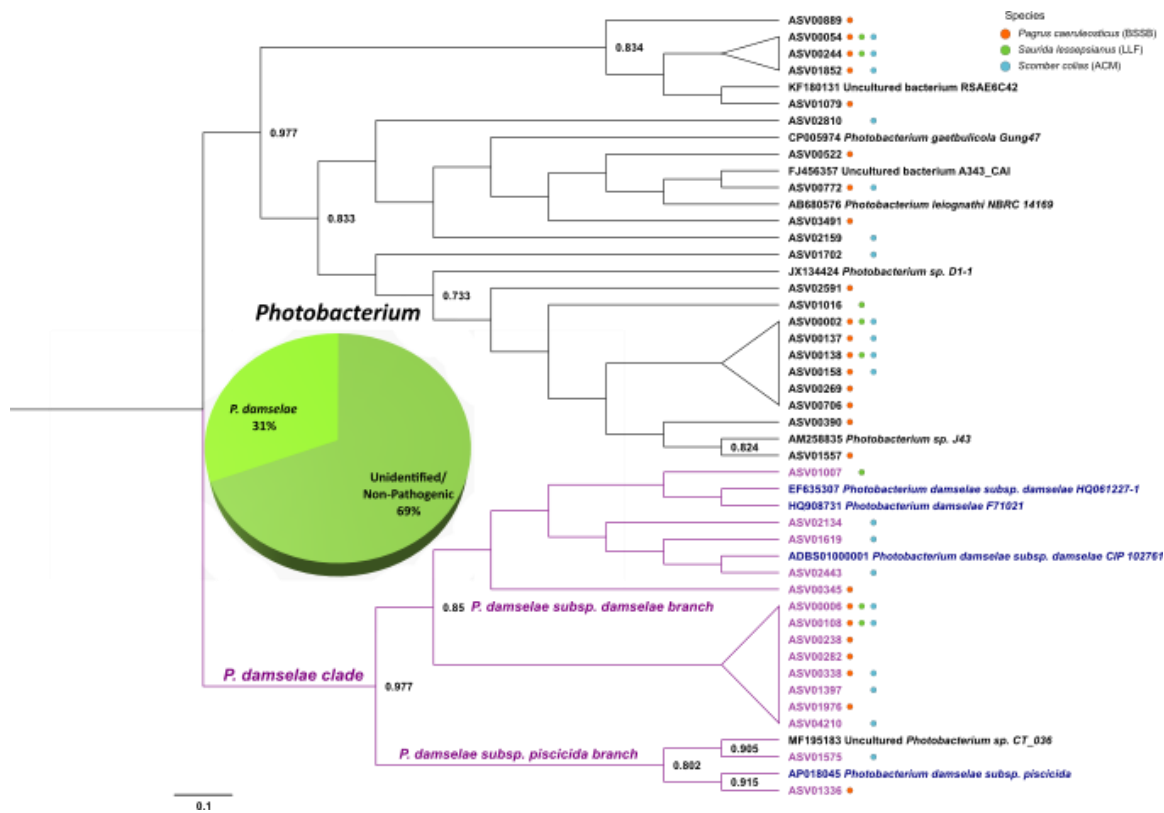


Figure S4: A phylogenetic tree for *Photobacterium*-related ASVs. A cutoff of 0.7 (70% bootstrap support) was made for nodes, thus any lower value is not presented. Triangular shaped tips represent sequences found to be practically identical. Grey-scale shapes (to the right of ASV numbers) represent fish species in which the ASVs appeared. Reference sequences include their GenBank accession numbers. Grey-colored sequence names mark pathogenic species and ASVs identified as bearing a similarity to pathogenic species-related sequences. The *P. damsela* clade branch and sub-branches are thickened, for emphasis. The pie chart refers to percentage out of total number of reads. Smaller ASV numbers indicate they were more common (in terms of total reads) than large number ASVs. The scale bar represents 0.1 nucleotide substitutions per site.

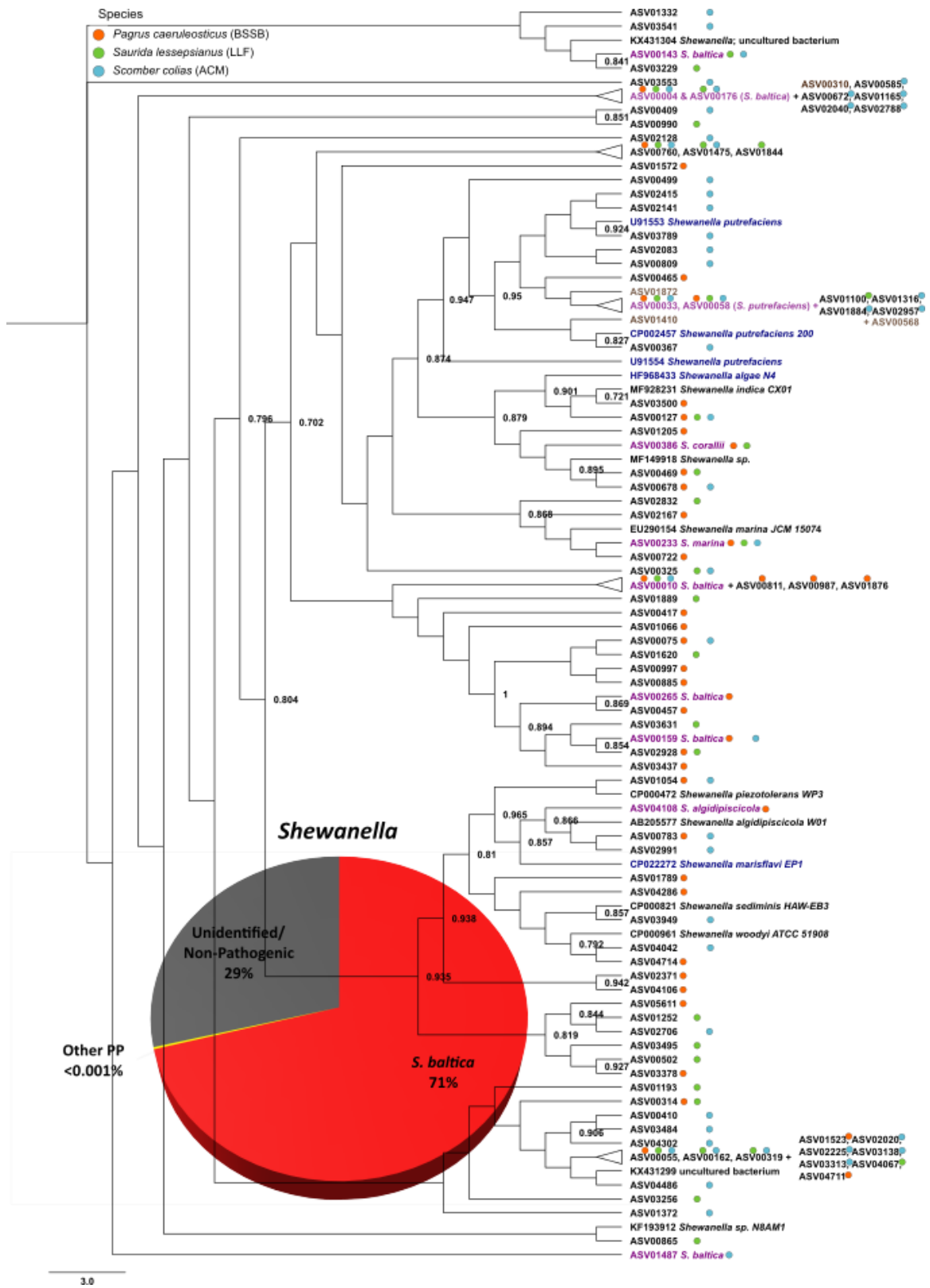


Figure S5: A phylogenetic tree for *Shewanella*-related ASVs. A cutoff of 0.7 (70% bootstrap support) was made for nodes, thus any lower value is not presented. Triangular shaped tips represent sequences found to be practically identical. Grey-scale shapes (above or to the right of ASV numbers) represent fish species in which the ASVs appeared. Reference sequences include their GenBank accession numbers. Grey-colored sequence names mark pathogenic species and ASVs identified as bearing a similarity to pathogenic species-related sequences, while light-grey ASVs mark those found in negative control samples only. The pie chart refers to percentage out of total number of reads. Smaller ASV numbers indicate they were more common (in terms of total reads) than large number ASVs. The scale bar represents 3.0 nucleotide substitutions per site.

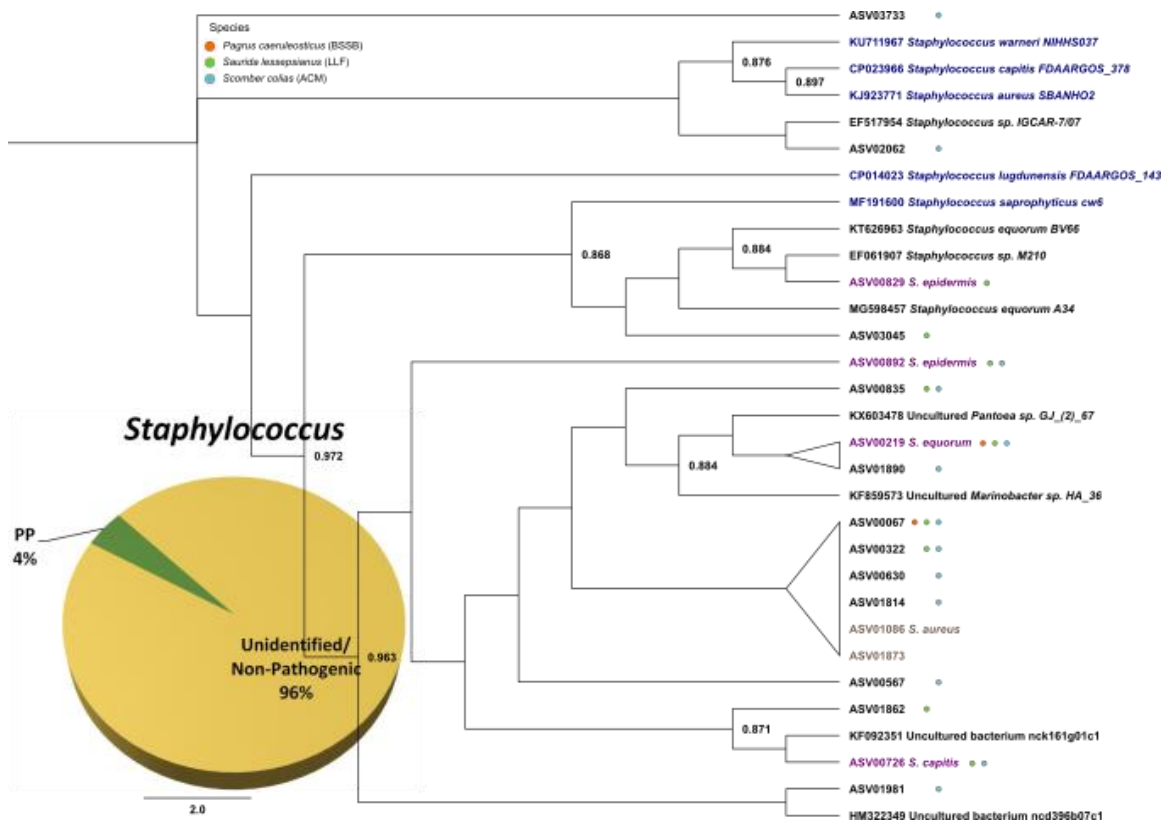


Figure S6: A phylogenetic tree for *Staphylococcus*-related ASVs. A cutoff of 0.7 (70% bootstrap support) was made for nodes, thus any lower value is not presented. Triangular shaped tips represent sequences found to be practically identical. Grey-scale shapes (to the right of ASV numbers) represent fish species in which the ASVs appeared. Reference sequences include their GenBank accession numbers. Grey-colored sequence names mark pathogenic species and ASVs identified as bearing a similarity to pathogenic species-related sequences, while light-grey ASVs mark those found in negative control samples only. The pie chart refers to percentage out of total number of reads. Smaller ASV numbers indicate they were more common (in terms of total reads) than large number ASVs. The scale bar represents 2.0 nucleotide substitutions per site.

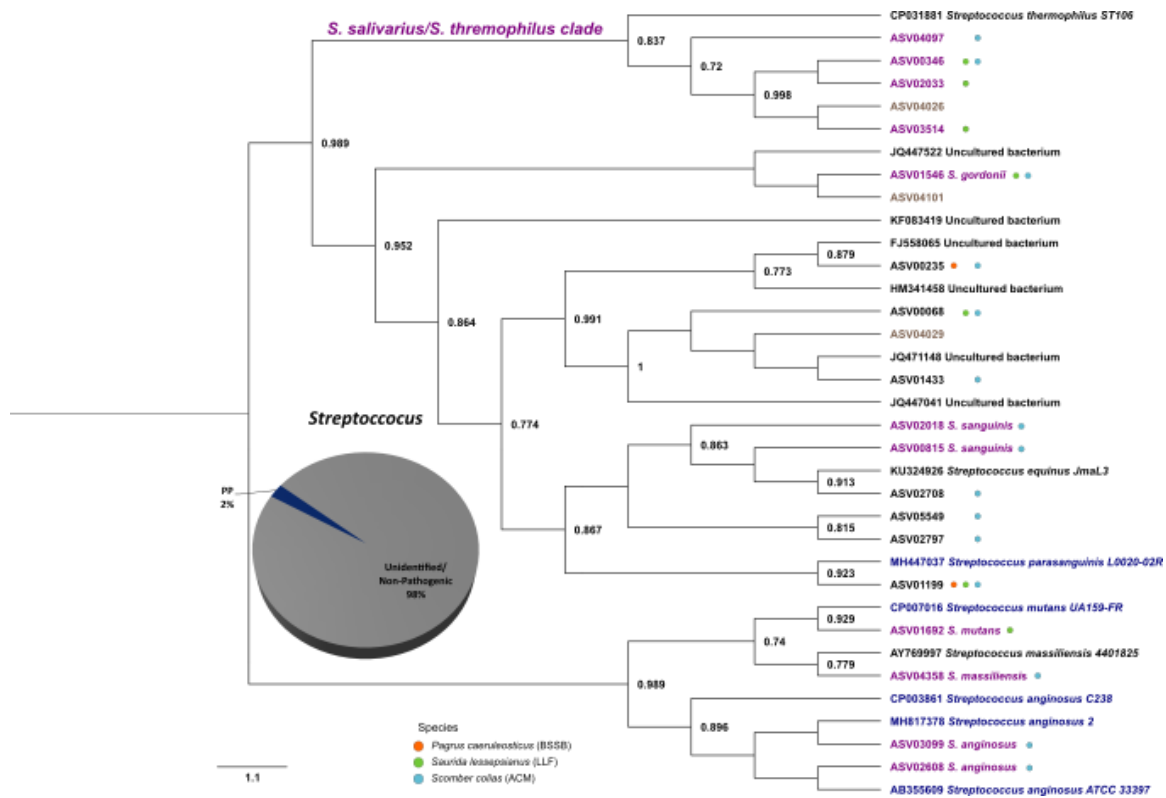


Figure S7: A phylogenetic tree for *Streptococcus*-related ASVs. A cutoff of 0.7 (70% bootstrap support) was made for nodes, thus any lower value is not presented. Triangular shaped tips represent sequences found to be practically identical. Grey-scale shapes (to the right of ASV numbers) represent fish species in which the ASVs appeared. Reference sequences include their GenBank accession numbers. Grey-colored sequence names mark pathogenic species and ASVs identified as bearing a similarity to pathogenic species-related sequences, while light-grey ASVs mark those found in negative control samples only. The pie chart refers to percentage out of total number of reads. Smaller ASV numbers indicate they were more common (in terms of total reads) than large number ASVs. The scale bar represents 1.1 nucleotide substitutions per site.

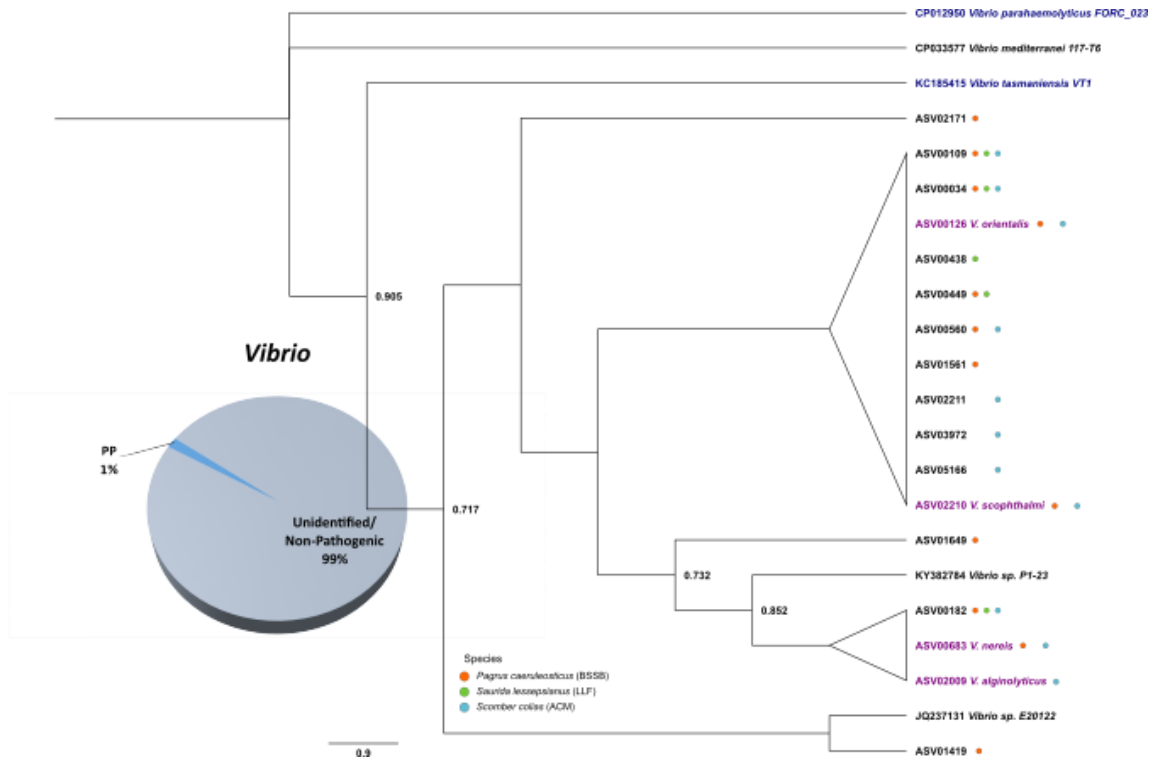


Figure S8: A phylogenetic tree for *Vibrio*-related ASVs. A cutoff of 0.7 (70% bootstrap support) was made for nodes, thus any lower value is not presented. Triangular shaped tips represent sequences found to be practically identical. Grey-scale shapes (to the right of ASV numbers) represent fish species in which the ASVs appeared. Reference sequences include their GenBank accession numbers. Grey-colored sequence names mark pathogenic species and ASVs identified as bearing a similarity to pathogenic species-related sequences. The pie chart refers to percentage out of total number of reads. Smaller ASV numbers indicate they were more common (in terms of total reads) than large number ASVs. The scale bar represents 0.9 nucleotide substitutions per site.

Supplementary Methods

Sequence data processing

Sequence data was analyzed using the Dada2 pipeline (72) using R package ‘dada2’ (version 1.14.1). Fastq formatted reads were trimmed and filtered for low quality using the command ‘filterAndTrim’ with parameters maxEE=2, maxN=0, trimleft=20 and the truncLen=150. Error rate estimation was carried out using the ‘learnError’ command with default parameters, but with the randomize parameter set to TRUE, in order to sample nucleotides and reads for model

building randomly across all samples. The dada2 algorithm was implemented for error correction and a count table containing the amplicon sequence variants and counts per sample was produced. Merging of forward and reverse reads was done using the 'mergePairs' command with a minimum overlap of 8 bases. Then, suspected chimeras were detected and removed using the command 'removeBimeraDenovo', with default parameters. For each amplicon sequence variant (ASV), taxonomy (up to the species level) was inferred by alignment to the Silva non-redundant small subunit ribosomal RNA database (version 138), using commands 'assignTaxonomy' and 'addSpecies' with default parameters, while setting the minimum bootstrap confidence value to 80%.

Data analysis

For data analysis and generation of figures, the online tool MicrobiomeAnalyst (<https://www.microbiomeanalyst.ca/MicrobiomeAnalyst/home.xhtml>) was used (73,74). Through Marker Data Profiling (MDP), three CSV files were uploaded – ASV counts, taxonomy, metadata. Taxonomy labels were assigned using the SILVA taxonomic framework (<https://www.arb-silva.de/documentation/silva-taxonomy/>). Data filtering settings were defined as follows: (i) minimum counts: 4; (ii) low count filter – prevalence in samples (%): 10; (iii) low variance filter – percentage to remove (%): 5, and filter based on inter-quantile range. This process removed 1093 low abundance/variance features and kept 283 others. For Alpha-diversity graphs data normalization settings were: (i) data rarefied to the minimum library size; (ii) no data scaling; and (iii) no data transformation. Alpha-diversity profiling settings were: (i) data input – filtered; (ii) experimental factor – 'species'; (iii) taxonomic level – 'feature-level'; (iv) diversity measure – observed, Shannon and Simpson; and (v) the statistical method – Mann-Whitney/Kruskal-Wallis. For all other graphs and analyses, the data normalization settings required not rarefying data, but transforming the data using relative log expression (RLE). For

taxa abundance stacked-bar graphs, graph type was set to percentage-abundance, merging small taxa with counts < 8,000 and grouping by fish species. For beta-diversity graphs, method: NMDS ordination, distance method calculation: Jaccard index. Statistical method: Permutational MANOVA (PERMANOVA), other options kept as default. For the correlation analysis: algorithm – Spearman's rank correlation, *P*-value threshold – 0.05, correlation threshold – 0.5. All other values used default settings.

Phylogenetic trees

Sequences identified as belonging to the several genera chosen for deeper enquiry were copied from the main data file converted to FASTA format, and then uploaded to Silva (<https://www.arb-silva.de/aligner/>) for preparing phylogenetic files (75–77). The ACT (Alignment, Classification and Tree Service) tool was used (SINA v1.2.11) (ref. 68), with the following parameters: (i) gene: ssu; (ii) unaligned remaining bases – attached to the last aligned base; (iii) search and classify (checked) with minimum identity with query sequence set to 0.98, and number of neighbors per query sequence at 2; (iv) compute tree (checked) with its workflow set to 'denovo including neighbors', 'FastTree' as the program to use, and 'gamma' as the rate model for likelihoods; (v) the output format: FASTA and file zip-compressed; and (vi) taxonomies selected for classification: SILVA and RDP (only). All other parameters were kept as default. Output TREE format files were extracted for visualization with the FigTree v1.4.3 software (<http://tree.bio.ed.ac.uk/software/figtree/>).